#### Support vector machines (SVMs) Lecture 2

### David Sontag New York University

Slides adapted from Luke Zettlemoyer, Vibhav Gogate, and Carlos Guestrin

# Geometry of linear separators (see blackboard)

A plane can be specified as the set of all points given by:



Vector from origin to a point in the plane

Two non-parallel directions in the plane



Alternatively, it can be specified as:  $(\mathbf{p} - \mathbf{a}) \cdot \mathbf{n} = 0 \Leftrightarrow \mathbf{p} \cdot \mathbf{n} = \mathbf{a} \cdot \mathbf{n}$ Normal vector (we will call this w)

Only need to specify this dot product, a scalar (we will call this the offset, b)

Barber, Section 29.1.1-4

## **Linear Separators**

- If training data is linearly separable, perceptron is guaranteed to find *some* linear separator
- Which of these is optimal?



#### Support Vector Machine (SVM)

SVMs (Vapnik, 1990's) choose the linear separator with the largest margin



- Good according to intuition, theory, practice
- SVM became famous when, using images as input, it gave accuracy comparable to neural-network with hand-designed features in a handwriting recognition task

Support vector machines: 3 key ideas

- 1. Use **optimization** to find solution (i.e. a hyperplane) with few errors
- 2. Seek **large margin** separator to improve generalization
- 3. Use **kernel trick** to make large feature spaces computationally efficient



# Finding a perfect classifier (when one exists) using linear programming



For every data point  $(x_t, y_t)$ , enforce the constraint

for y<sub>t</sub> = +1, 
$$w \cdot x_t + b \ge 1$$
  
and for y<sub>t</sub> = -1,  $w \cdot x_t + b \le -1$ 

Equivalently, we want to satisfy all of the linear constraints

$$y_t \left( w \cdot x_t + b \right) \ge 1 \quad \forall t$$

This *linear program* can be efficiently solved using algorithms such as simplex, interior point, or ellipsoid

# Finding a perfect classifier (when one exists) using linear programming

Example of 2-dimensional linear programming (feasibility) problem:

For SVMs, each data point gives one inequality:  $y_t (w \cdot x_t + b) \ge 1$ 



#### What happens if the data set is not linearly separable?

#### Minimizing number of errors (0-1 loss)



• Try to find weights that violate as few constraints as possible?

 $\begin{array}{ll} \text{minimize}_{\mathbf{w},b} & \text{#(mistakes)} \\ \left(\mathbf{w}.\mathbf{x}_{j}+b\right)y_{j} \geq 1 & ,\forall j \end{array}$ 

• Formalize this using the 0-1 loss:

$$\begin{split} \min_{\mathbf{w},b} \sum_{j} \ell_{0,1}(y_j, \, w \cdot x_j + b) \\ \text{where } \ell_{0,1}(y, \hat{y}) = \mathbf{1}[y \neq \operatorname{sign}(\hat{y})] \end{split}$$

- Unfortunately, minimizing 0-1 loss is NP-hard in the worst-case
  - Non-starter. We need another approach.

#### Key idea #1: Allow for slack



#### For each data point:

- •If functional margin  $\geq$  1, don't care
- •If functional margin < 1, pay linear penalty

#### Key idea #1: Allow for *slack*



#### Equivalent hinge loss formulation

$$egin{aligned} \mathsf{minimize}_{\mathbf{w},b,\,\xi} & \Sigma_{\mathrm{j}}\,\xi_{\mathrm{j}} \ & \left(\mathbf{w}.\mathbf{x}_{j}+b
ight)y_{j} \geq \mathbf{1}$$
 -  $\xi_{\mathrm{j}}$  ,  $orall j$   $\xi_{\mathrm{j}} \geq 0$ 

Substituting  $\xi_j = \max(0, 1 - (w \cdot x_j + b) y_j)$  into the objective, we get:

$$\min_{w,b} \sum_{j} \max\left(0, 1 - \left(w \cdot x_{j} + b\right) y_{j}\right)$$

The hinge loss is defined as  $\ell_{\text{hinge}}(y, \hat{y}) = \max\left(0, 1 - \hat{y}y\right)$ 

$$\min_{\mathbf{w},b} \sum_{j} \ell_{\text{hinge}}(y_j, \, w \cdot x_j + b)$$

This is empirical risk minimization, using the hinge loss

#### Hinge loss vs. 0/1 loss



#### Hinge loss upper bounds 0/1 loss!

It is the tightest *convex* upper bound on the 0/1 loss

#### Key idea #2: seek large margin



#### Key idea #2: seek large margin

 Suppose again that the data is linearly separable and we are solving a feasibility problem, with constraints

$$y_t \left( w \cdot x_t + b \right) \ge 1 \quad \forall t$$

 If the length of the weight vector ||w|| is too small, the optimization problem is infeasible! Why?



#### What is $\gamma$ (geometric margin) as a function of w?



Final result: can maximize  $\gamma$  by minimizing  $||w||_2!!!$ 

#### (Hard margin) support vector machines



 $\begin{array}{ll} \text{minimize}_{\mathbf{w},b} & \mathbf{w}.\mathbf{w} \\ \left(\mathbf{w}.\mathbf{x}_{j}+b\right)y_{j} \geq \mathbf{1}, \ \forall j \end{array}$ 

#### Example of a **convex optimization** problem

- A quadratic program
- Polynomial-time algorithms to solve!
- Hyperplane defined by support vectors
  - Could use them as a lower-dimension basis to write down line, although we haven't seen how yet

#### More on these later

#### Allowing for slack: "Soft margin SVM"



$$\begin{array}{ll} \min i z \mathbf{e}_{\mathbf{w},b} & \mathbf{w}.\mathbf{w} + C \Sigma_{j} \xi_{j} \\ \left(\mathbf{w}.\mathbf{x}_{j} + b\right) y_{j} \geq \mathbf{1} - \xi_{j} &, \forall j \xi_{j} \geq 0 \\ & \uparrow \\ & \text{``slack variables''} \end{array}$$

#### Slack penalty C > 0:

- $C=\infty \rightarrow$  have to separate the data!
- $C=0 \rightarrow$  ignores the data entirely!
- Select using cross-validation

For each data point:

- •If margin  $\geq$  1, don't care
- •If margin < 1, pay linear penalty

#### Equivalent formulation using hinge loss

$$\begin{array}{ll} \text{minimize}_{\mathbf{w},b} & \mathbf{w}.\mathbf{w} + C \Sigma_{j} \xi_{j} \\ \left(\mathbf{w}.\mathbf{x}_{j} + b\right) y_{j} \geq \mathbf{1} - \xi_{j} &, \forall j \ \xi_{j} \geq 0 \end{array}$$

Substituting  $\xi_j = \max(0, 1 - (w \cdot x_j + b) y_j)$  into the objective, we get:

$$\min ||w||^2 + C \sum_j \max (0, 1 - (w \cdot x_j + b) y_j)$$

The hinge loss is defined as  $\ell_{\text{hinge}}(y, \hat{y}) = \max\left(0, 1 - \hat{y}y\right)$ 

$$\min_{\mathbf{w},b} ||w||_2^2 + C \sum_j \ell_{\text{hinge}}(y_j, w \cdot x_j + b)$$

This is called **regularization**; used to prevent overfitting!

This part is empirical risk minimization, using the hinge loss

# What if the data is not linearly separable?

Use features of features of features....



$$\phi(x) = \begin{pmatrix} x^{(1)} & \\ & \ddots & \\ & x^{(n)} \\ x^{(1)}x^{(2)} \\ & x^{(1)}x^{(3)} \\ & \ddots & \\ & e^{x^{(1)}} \\ & \ddots & \end{pmatrix}$$

Feature space can get really large really quickly!





Non-linear separator in the original x-space



Linear separator in the feature  $\phi$ -space

[Tommi Jaakkola]

## What's Next!

- Learn one of the most interesting and exciting recent advancements in machine learning
  - Key idea #3: the "kernel trick"
  - High dimensional feature spaces at no extra cost
- But first, a detour
  - Constrained optimization!



How do we solve with constraints? → Lagrange Multipliers!!!

#### Lagrange multipliers – Dual variables





We will solve:

 $\min_x \max_\alpha L(x, \alpha)$ 

#### Why is this equivalent?

• min is fighting max! x<b  $\rightarrow$  (x-b)<0  $\rightarrow$  max<sub> $\alpha$ </sub>- $\alpha$  (x-b) =  $\infty$ 

min won't let this happen!

 $(x-b) = \infty$  S.t.  $\alpha \ge 0$ appen!

Add new constraint

x>b,  $\alpha \ge 0 \rightarrow (x-b) > 0 \rightarrow \max_{\alpha} - \alpha (x-b) = 0$ ,  $\alpha *=0$ 

• min is cool with 0, and L(x,  $\alpha$ )=x<sup>2</sup> (original objective)

 $x=b \rightarrow \alpha$  can be anything, and L(x,  $\alpha$ )=x<sup>2</sup> (original objective)

The *min* on the outside forces *max* to behave, so constraints will be satisfied.

Dual SVM derivation (1) – the linearly separable case (hard margin SVM)

**Original optimization problem:** 



Our goal now is to solve:  $\min_{\vec{w},b} \max_{\vec{\alpha} \ge 0} L(\vec{w},\vec{\alpha})$ 

Dual SVM derivation (2) – the linearly separable case (hard margin SVM)

(Primal) 
$$\min_{\vec{w},b} \max_{\vec{\alpha} \ge 0} \frac{1}{2} ||\vec{w}||^2 - \sum_j \alpha_j \left[ (\vec{w} \cdot \vec{x}_j + b) \, y_j - 1 \right]$$
  
Swap min and max  
(Dual) 
$$\max_{\vec{\alpha} \ge 0} \min_{\vec{w},b} \frac{1}{2} ||\vec{w}||^2 - \sum_j \alpha_j \left[ (\vec{w} \cdot \vec{x}_j + b) \, y_j - 1 \right]$$

*Slater's condition* from convex optimization guarantees that these two optimization problems are equivalent!

Dual SVM derivation (3) – the linearly separable case (hard margin SVM)

(Dual) 
$$\max_{\vec{\alpha} \ge 0} \min_{\vec{w}, b} \frac{1}{2} ||\vec{w}||^2 - \sum_j \alpha_j \left[ (\vec{w} \cdot \vec{x}_j + b) y_j - 1 \right]$$

Can solve for optimal **w**, b as function of  $\alpha$ :

$$\frac{\partial L}{\partial w} = w - \sum_{j} \alpha_{j} y_{j} x_{j} \quad \Rightarrow \quad \mathbf{w} = \sum_{j} \alpha_{j} y_{j} \mathbf{x}_{j}$$
$$\frac{\partial L}{\partial b} = -\sum_{j} \alpha_{j} y_{j} \quad \Rightarrow \quad \sum_{j} \alpha_{j} y_{j} = 0$$

Substituting these values back in (and simplifying), we obtain:

(Dual) 
$$\vec{\alpha} \ge 0, \sum_{j} \alpha_{j} y_{j} = 0$$
  $\sum_{j} \alpha_{j} - \frac{1}{2} \sum_{i,j} y_{i} y_{j} \alpha_{i} \alpha_{j} (\vec{x}_{i} \cdot \vec{x}_{j})$   
Sums over all training examples scalars dot product

Dual formulation only depends on dot-products of the features!

$$\max_{\vec{\alpha} \ge 0, \sum_{j} \alpha_{j} y_{j} = 0} \quad \sum_{j} \alpha_{j} - \frac{1}{2} \sum_{i,j} y_{i} y_{j} \alpha_{i} \alpha_{j} \left( \vec{x}_{i} \cdot \vec{x}_{j} \right)$$

First, we introduce a *feature mapping*:

$$\mathbf{x}_i \mathbf{x}_j \rightarrow \Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j)$$

Next, replace the dot product with an equivalent kernel function:

maximize<sub>$$\alpha$$</sub>  $\sum_{i} \alpha_{i} - \frac{1}{2} \sum_{i,j} \alpha_{i} \alpha_{j} y_{i} y_{j} K(\mathbf{x}_{i}, \mathbf{x}_{j})$   
 $K(\mathbf{x}_{i}, \mathbf{x}_{j}) = \Phi(\mathbf{x}_{i}) \cdot \Phi(\mathbf{x}_{j})$   
 $\sum_{i} \alpha_{i} y_{i} = 0$ 

#### SVM with kernels

maximize<sub>$$\alpha$$</sub>  $\sum_{i} \alpha_{i} - \frac{1}{2} \sum_{i,j} \alpha_{i} \alpha_{j} y_{i} y_{j} K(\mathbf{x}_{i}, \mathbf{x}_{j})$   
 $K(\mathbf{x}_{i}, \mathbf{x}_{j}) = \Phi(\mathbf{x}_{i}) \cdot \Phi(\mathbf{x}_{j})$   
 $\sum_{i} \alpha_{i} y_{i} = 0$   
 $C > \alpha_{i} > 0$ 

- Never compute features explicitly!!!
  - Compute dot products in closed form

Predict with:

$$y \leftarrow \operatorname{sign}\left[\sum_{i} \alpha_{i} y_{i} K(x_{i}, x) + b\right]$$

- O(n<sup>2</sup>) time in size of dataset to compute objective
  - much work on speeding up

## Common kernels

- Polynomials of degree exactly d $K(\mathbf{u}, \mathbf{v}) = (\mathbf{u} \cdot \mathbf{v})^d$
- Polynomials of degree up to *d*

$$K(\mathbf{u},\mathbf{v}) = (\mathbf{u} \cdot \mathbf{v} + 1)^d$$

Gaussian kernels

$$K(\vec{u}, \vec{v}) = \exp\left(-\frac{||\vec{u} - \vec{v}||_2^2}{2\sigma^2}\right)$$

• Sigmoid

$$K(\mathbf{u},\mathbf{v}) = \tanh(\eta \mathbf{u} \cdot \mathbf{v} + \nu)$$

• And many others: very active area of research!

#### **Quadratic kernel**



Non-linear separator in the original x-space



Linear separator in the feature  $\phi$ -space

[Tommi Jaakkola]

#### **Quadratic kernel**

$$k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z} + c)^2 = \left(\sum_{j=1}^n x^{(j)} z^{(j)} + c\right) \left(\sum_{\ell=1}^n x^{(\ell)} z^{(\ell)} + c\right)$$
$$= \sum_{j=1}^n \sum_{\ell=1}^n x^{(j)} x^{(\ell)} z^{(j)} z^{(\ell)} + 2c \sum_{j=1}^n x^{(j)} z^{(j)} + c^2$$
$$= \sum_{j,\ell=1}^n (x^{(j)} x^{(\ell)}) (z^{(j)} z^{(\ell)}) + \sum_{j=1}^n (\sqrt{2c} x^{(j)}) (\sqrt{2c} z^{(j)}) + c^2,$$

Feature mapping given by:

$$\mathbf{\Phi}(\mathbf{x}) = [x^{(1)2}, x^{(1)}x^{(2)}, ..., x^{(3)2}, \sqrt{2c}x^{(1)}, \sqrt{2c}x^{(2)}, \sqrt{2c}x^{(3)}, c]$$

[Cynthia Rudin]